# Ziggy, a Portable, Scalable Infrastructure for Science Data Processing Pipelines and its Application to a Proxy, Legacy Global Hyperspectral Data Set for NASA's Earth System Observatory's Upcoming Surface, Biology and Geology Mission

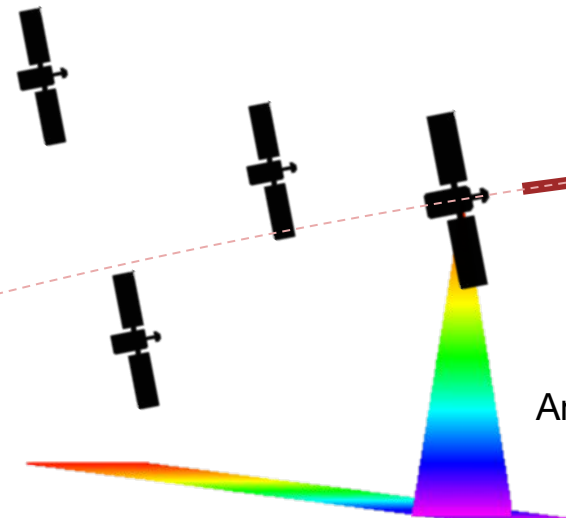Jon M. Jenkins*

Advanced Supercomputing Division

NASA Ames Research Center

24 March 2023

IEEE SPIN 2023

Amity University, New Delhi, India and Everywhere

*With much help from the Ames SBG/SISTER Team

National Aeronautics and Space Administration

# Overview

- Science Data Systems are Expensive – is there a better way?
- Building science pipelines for astrophysics missions
- Ziggy, a flexible scalable infrastructure for science pipelines
- An overview of the Surface Biology & Geology (SBG) Mission
- Resurrecting the Hyperion/Earth Observer-1 science workflow
- Conclusions

# Designing, Developing and Deploying New Science Data Systems is Expensive

- **Case Study:** Kepler → TESS
  - The total cost of the first complete, operational version of the Kepler SDS was 78 person years, or ~$25M against a $300M cost cap
  - The cost of the pipeline infrastructure was 40 person years or ~$13M of this cost
  - We reused ~60% of the Kepler codebase to develop the TESS codebase
  - This reduced the cost of the TESS pipeline to ~31 person years or about $10M
- **Can we leverage our Kepler/TESS experience to be able to build complex science pipelines for much less cost and schedule risk? Yes!**
- **The Solution: Ziggy, an open-source flexible, scalable infrastructure for science pipelines**
  - We extracted the pipeline infrastructure components of the TESS pipeline and set up a new software project for Ziggy
  - It only took 1 week to build a simple demo pipeline for AVIRIS data with an open-source science algorithm in Python
  - We've spent the last ~2.5 years building additional functionality into Ziggy with a small team of 2 people
  - Now we can build new science pipelines in a matter of days to weeks instead of months to years
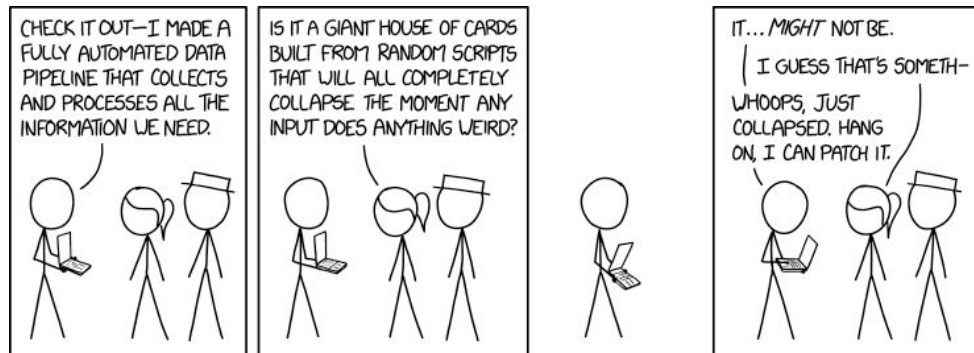
# Kepler Science Operations Center Architecture

The Kepler SOC Pipeline consisted of 24 Computer Software Configuration Items, including the Pipeline Infrastructure (PI), which has been split off as its own software project, *Ziggy*

*Ziggy* is capable of running complex science algorithms on large data sets on a laptop, workstation or on the NAS Pleiades supercomputer.

We've demonstrated running Ziggy in AWS

TESS SOC directly connected to NAS Pleaides backbone network

Performant Access to Compute and Filesystems

# ziggy — A Flexible, Scalable Pipeline Controller

Highly automated processing of large volumes of data

*github.com/nasa/ziggy

Automated dispatching to Pleiades for large to extremely large tasks

Provides permanent storage for all version changes, linkage between parameter values and pipeline tasks

Excellent logging and diagnostics for those rare occasions when something goes wrong

Extensively tested by Kepler and TESS processing missions

Codebase: Java (some C, C++, Perl, Shell scripts)

Relational database: mainly for data accountability

Datastore: main storage of mission data and products

Pipeline definitions: XML files that sequence algorithms, specify parameters, etc.

Messaging: Remote Method Invocation (a Java API)

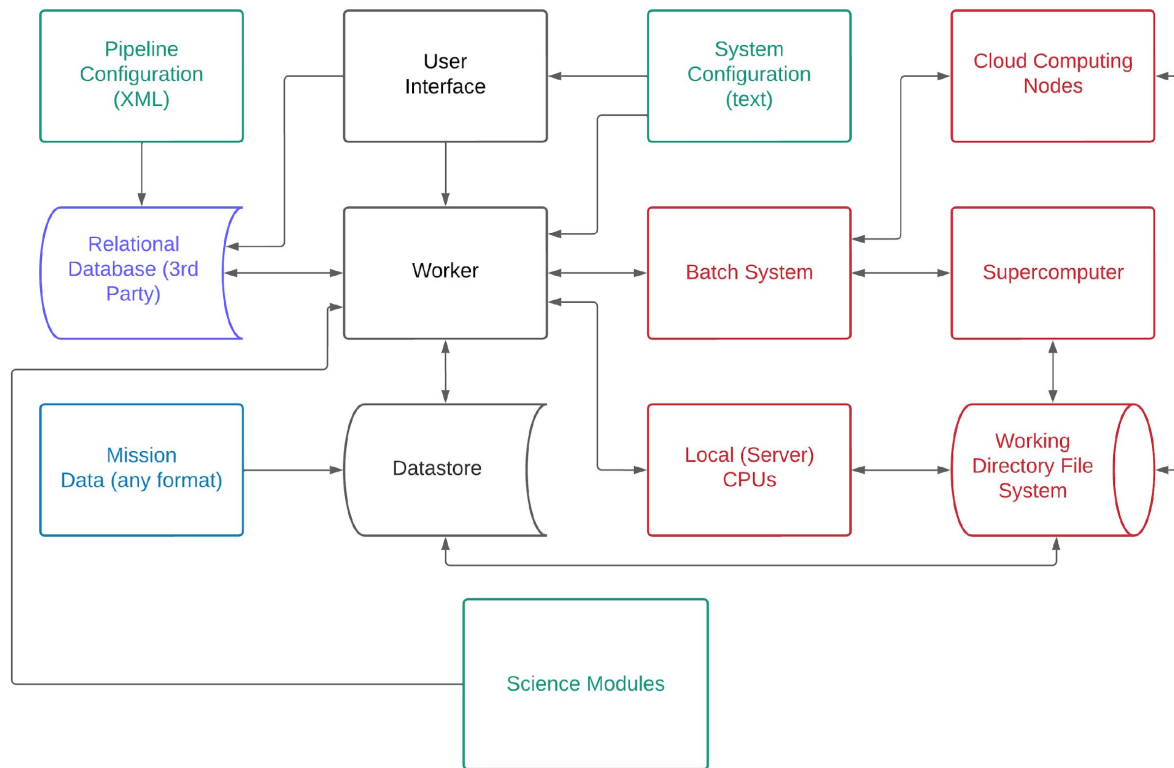Ziggy is approved for release under the NASA Open Source Software Initiative



https://xkcd.com/2054/ This work is licensed under a Creative Commons Attribution-NonCommercial 2.5 License

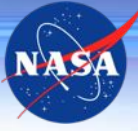Ziggy is TRL 7 and a class C software under NPR 7150.2C

# Architecture



Ziggy scales from laptops to HPC and to Cloud
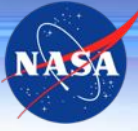
**ziggy**



**Develop here…**

**… run here!**

Ziggy is sufficiently lightweight to run on a laptop and sufficiently robust to run on a supercomputer; builds on Mac OS X and Linux are supported.

We've built and operated several major science pipelines for Astrophysics missions.



Credit: S. Quinn

We've built and operated several major science pipelines for Astrophysics missions.

Kepler Mission:
Discovered 3251 of the 5300 known exoplanets to date
1 GB day$^{-1}$



*As per NExScI data base as of 9/14/22

National Aeronautics and Space Administration
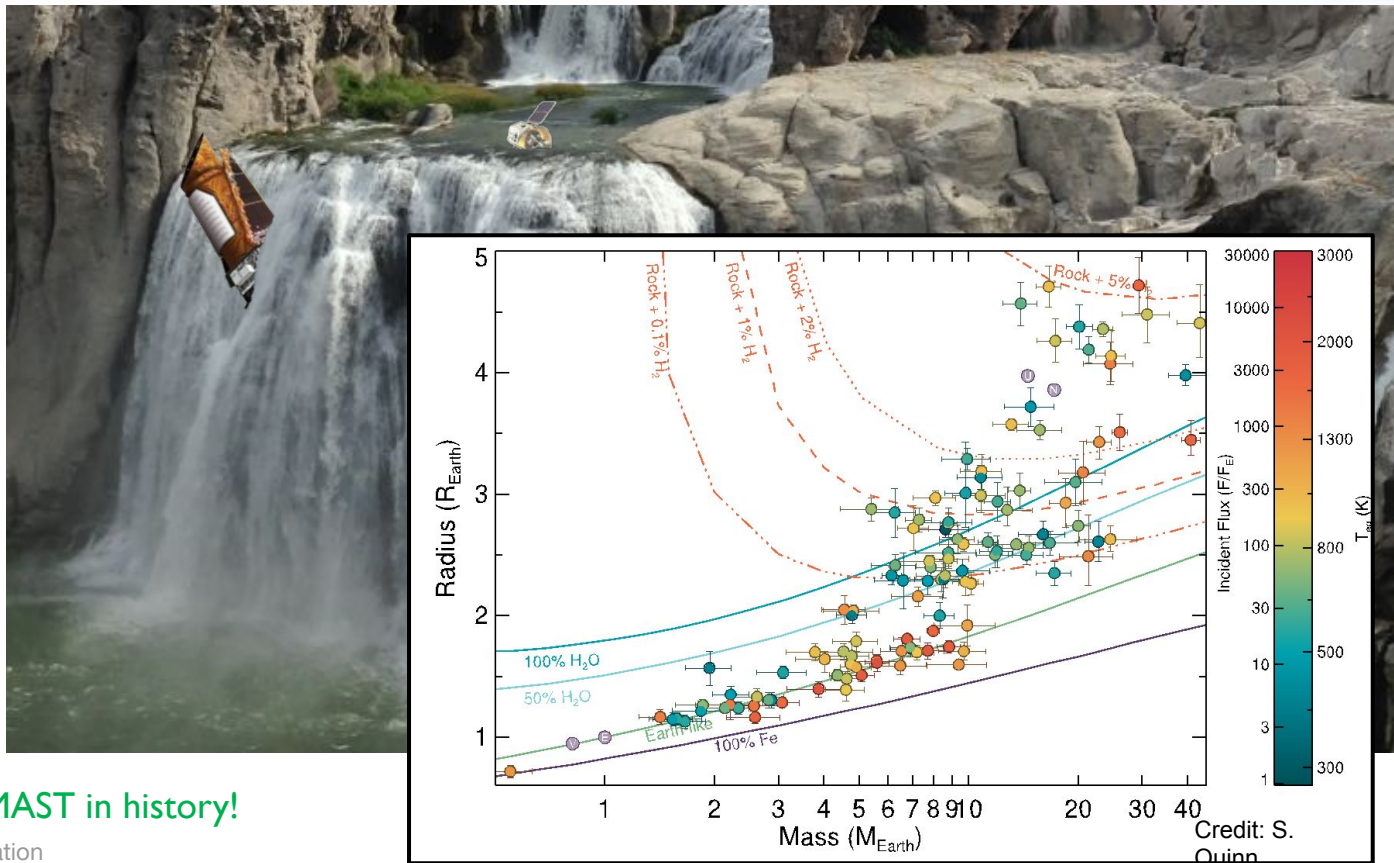
https://exoplanetarchive.ipac.caltech.edu/

We've built and operated several major science pipelines for Astrophysics missions.

Kepler Mission:
Discovered 3251 of the 5300 known exoplanets to date
1 GB day$^{-1}$

Transiting Exoplanet Survey Satellite:
316 Planets
104 < 2.5 $R_{earth}$ with measured masses
13-50 GB day$^{-1}$

Most downloaded data set at MAST in history!

Credit: S. Quinn

## SBG is key to understanding in five research and applications focus areas:

- Terrestrial and aquatic ecosystems
- Hydrology
- Weather
- Climate
- Solid Earth

## The Decadal Survey defines the implementation as two sensors *"Hyperspectral imagery in the visible and shortwave infrared; multi- or hyperspectral imagery in the thermal IR"*:

1. "... .a moderate spatial resolution (30-45 m GSD), hyperspectral resolution (10 nm; 400-2500 nm), high fidelity (SNR = 400:1 VNIR/250:1 SWIR) imaging spectrometer is needed for characterizing land, inland aquatic, coastal zone, and shallow coral reef ecosystems"
2. "... .30-60 m TIR observations in the 10.5-11.5 μm and 11.5-12.5 μm spectral regions are needed with a 2-4 day revisit frequency"[1]

1) Note, this specification was updated based on recent work and community engagement to optimize for the DS-specified science and applications.

# SBG Architecture

**SBG Constellation Pathfinder**

agenzia spaziale italiana

SBG's per-day volume will be greater than NASA's total extant airborne hyperspectral data collection!

**SBG Light**
Wide-swath VSWIR spectrometer

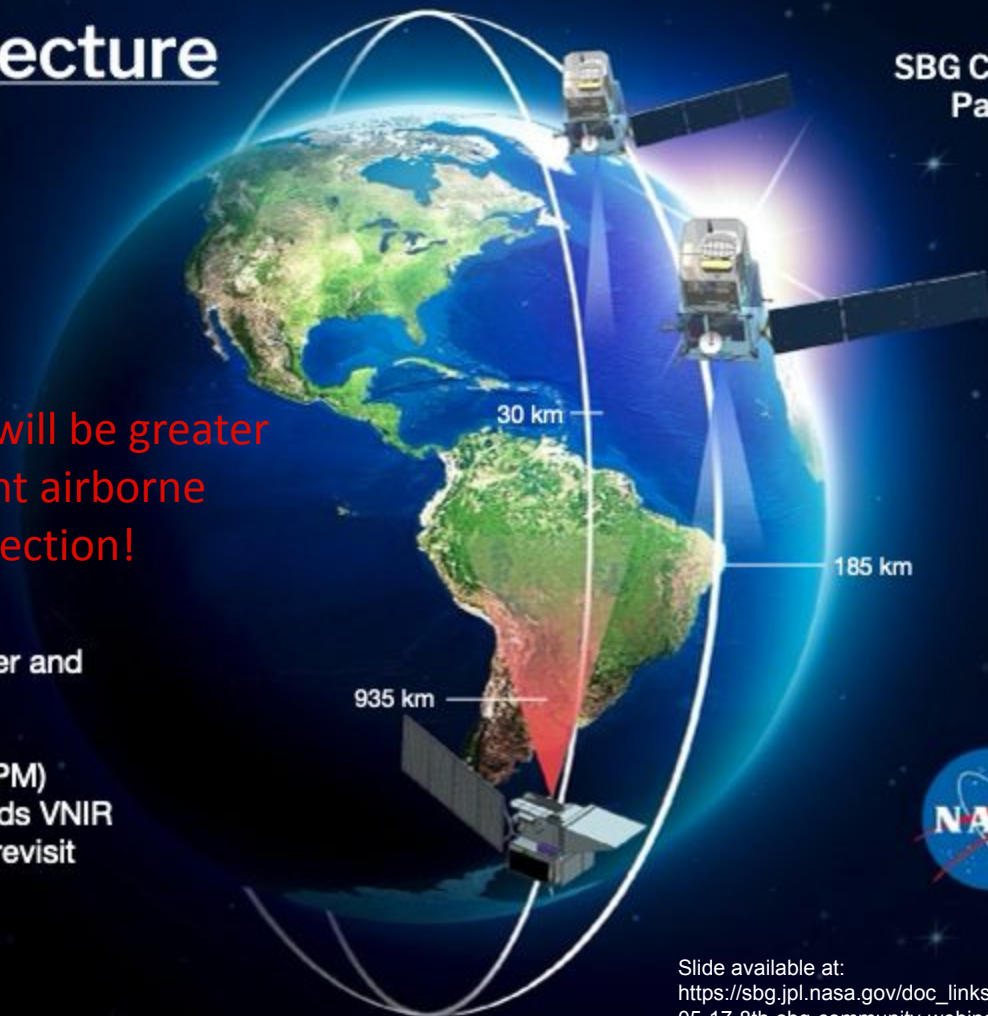Sun-sync orbit (late AM)
185 km swath
16 day revisit
10 nm, 200+ bands
30 meter GSD
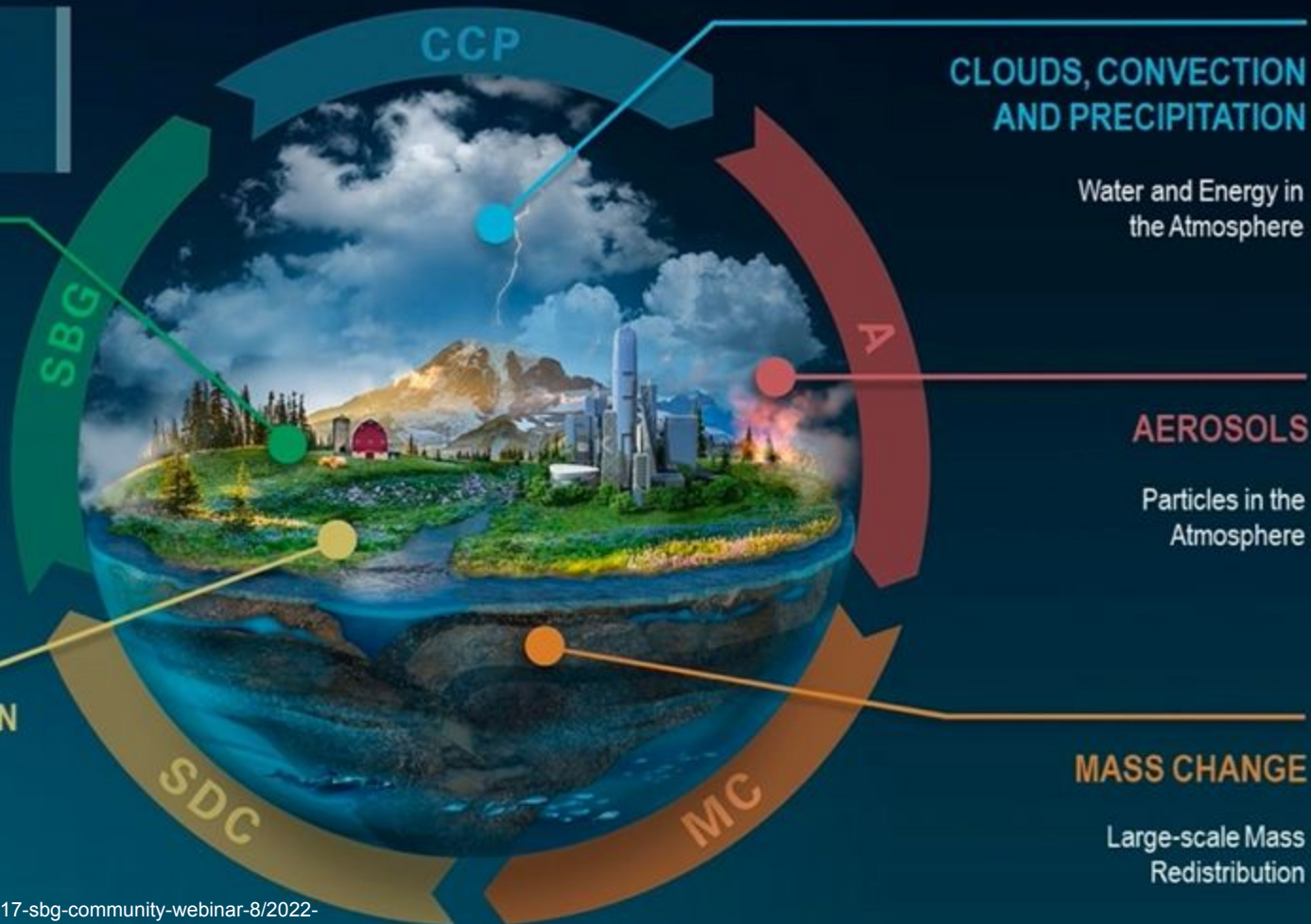High SNR and radiometric performance
~5 deg off-nadir tilt

**SBG Heat**
Wide-swath TIR imager and ASI VNIR camera

Sun-sync orbit (early PM)
5+ bands TIR, 2+ bands VNIR
935 km swath, 3 day revisit
60 meter GSD
0.2K NeDT

30 km
185 km
935 km

NASA

2

EARTH SYSTEM OBSERVATORY

CCP

CLOUDS, CONVECTION AND PRECIPITATION

Water and Energy in the Atmosphere

SURFACE BIOLOGY AND GEOLOGY

Earth Surface & Ecosystems

SBG

A

AEROSOLS

Particles in the Atmosphere

SURFACE DEFORMATION AND CHANGE

Earth Surface Dynamics

SDC

MC

MASS CHANGE

Large-scale Mass Redistribution
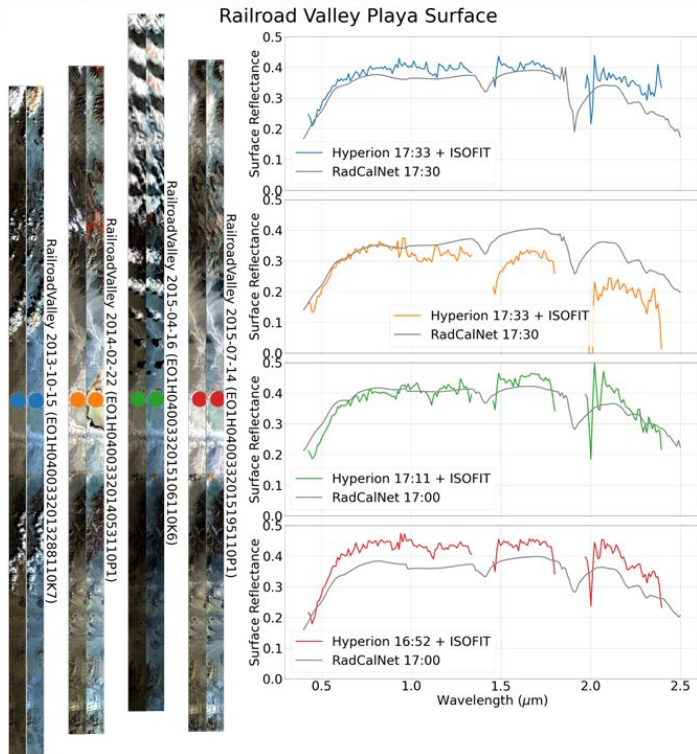
The Space-based Imaging Spectroscopy and Thermal pathfindER (SISTER*) study seeks to adapt and refine existing workflows and architectures as a prototype for the SBG science data system to reduce technical risk.

Hyperion Pipeline:
- We've resurrected the 17-year, 55-TB Hyperion pipeline and reprocessed nearly the entire data set to L2 (surface reflectance)
- Currently checking consistency of Hyperion surface reflectance results (L2) (RadCalNet and AVIRIS/Hyperion Comparisons)
- Future work:
  – Incorporate ground control database and georeferencing and co-registration
  – Incorporate L3 algorithms for vegetative traits and/or aquatic studies
- Note that SBG will collect 2.4 TB day-1 of hyperspectral image data and will produce ~40 TB day-1 of data products



**Hyperion Processing Flow**
June 23, 2022

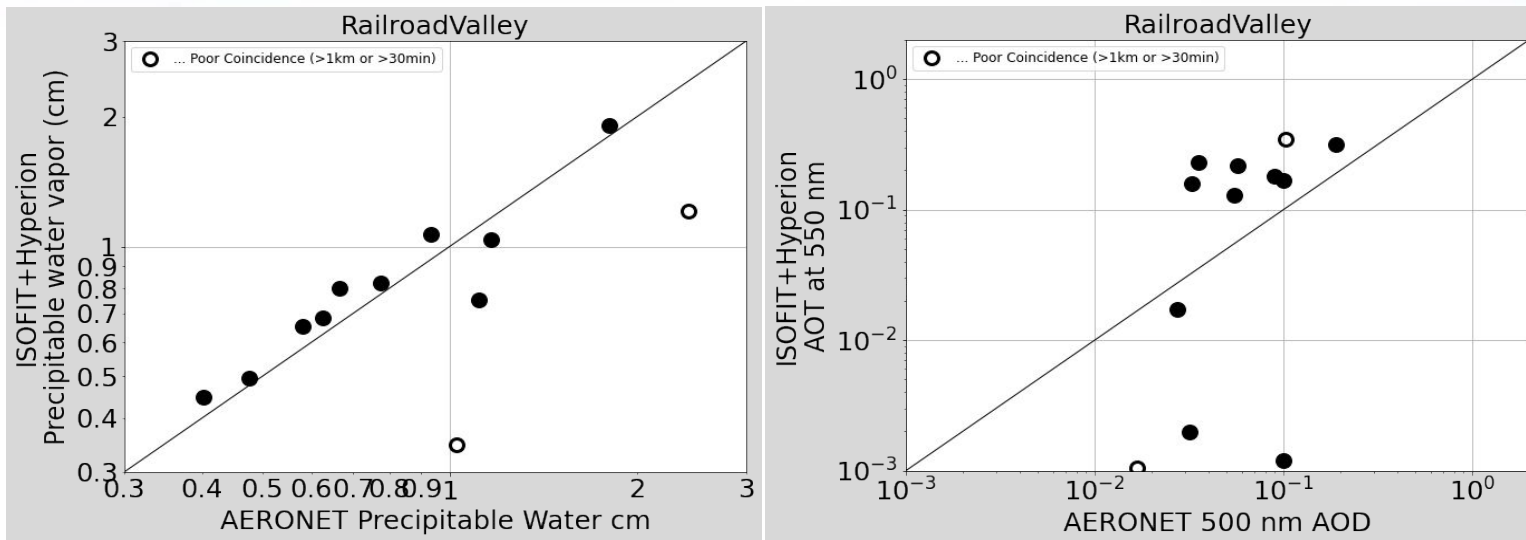Comparison of RadCalNet measurements with Hyperion surface reflectance retrievals for scenes observed in Railroad Valley.



Comparison of surface reflectance spectra retrieved with the Hyperion and AVIRIS sensors, each with ISOFIT and ATREM retrieval algorithms. The results are shown for a vegetation site near Half Moon Bay, CA, observed on April 30, 2015.
Expected discrepancies annotated according to Petya Campbell.

Comparison of Hyperion ISOFIT atmospheric products (left: water vapor, right: aerosol optical depth) with AERONET observations at Railroad Valley.

The RMS differences for near-coincident cases (< 1 km horizontal distance and <30-minute temporal gap, including scenes not shown in Figure 2) marked with closed circles in Figure 3 is 0.14 cm in column water vapor and 0.11 in 500-550 nm aerosol optical depth, respectively.

The aerosol optical depth retrieved with Hyperion and ISOFIT tends to be greater than the ground-observed value. Systematic differences may be explained by calibration inaccuracy; their dependence on wavelength, location and year remain to be investigated.

# HECC Forward and Reprocessing Cost Est

- Costs Basis Est from last customer to buy HPC on Margin

  $0.10/SBU - 7 year

- Scaling Hyperion up to SBG

- 55 TB L0 ⇔ ~22 days of SBG

- 14000 SBUs for Hyperion is ~232,432 SBUs/year for SBG

- 7-year purchase, 64 Milan nodes, 2.5 PB storage

- 2 PB of tape each year for L0 (2 copies)

- Reprocessing every two years

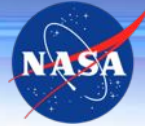- 2 hours compute/2 hours transfer @30gb/sec

- 2 hours compute/2 hours transfer 69% of capacity available: algorithm development, >L2 processing, HySDS hybrid post DAAC analysis, more frequent re-processing, open science activities, etc.

| | Forward (KiloSBU) | Reprocessing (KiloSBU) | Annual (KiloSBU) | Compute Tape | Data Storage (K $) | Projected HECC Cost (K $) |
|---|---|---|---|---|---|---|
| Year | | | | | | |
| 0 | | | | $ 1,719 | $ 250 | $ 1,969 |
| 1 | 232 | | 232 | $ 84 | | $ 84 |
| 2 | 232 | | 232 | $ 84 | | $ 84 |
| 3 | 232 | 464 | 696 | $ 84 | | $ 84 |
| 4 | 232 | | 232 | $ 84 | | $ 84 |
| 5 | 232 | 928 | 1160 | $ 84 | | $ 84 |
| 6 | 232 | | 232 | $ 84 | $ 25 | $ 109 |
| 7 | 232 | 1392 | 1856 | $ 84 | $ 25 | $ 109 |
| | | | | | | |
| | | | | | | $ 2,607 |
| | | | | | percent utilization | 31% |

# Conclusions

- Our open-source science pipeline, Ziggy, is a key enabling technology, greatly reducing the cost, risk and schedule for developing and deploying new complex science pipelines

- The Surface, Biology and Geology Mission will revolutionize our understanding of the Earth that benefit humankind in many ways:

  - Agriculture, food security and surface water management
  - Water quality and coastal zones
  - Conservation
  - Wildfire risk and recovery
  - Disasters and natural hazards
  - Geology applications

- We've resurrected the Hyperion pipeline and reprocessed the 55-TB data set to top-of-the-atmosphere radiance and a selected portion to surface reflectance spectra

- The results of the Hyperion pathfinder pipeline study indicate that using Ziggy, NASA's Advanced Supercomputing Division can scale to support SBG's challenging data rate and computation requirements at a relatively low cost to NASA compared to commercial cloud options

# Backup Slides

# HPC Cost Terminology vs Cloud Costs

- **SBU\* – Standard Billing Unit**

  - SBU – Standard Billing Unit – Work completed in 1 hour on a dual socket Broadwell Node
  - SBU Cost == 75% of all HECC Compute capacity available day 1 of FY / FY Projected Budget
  - $0.56 in FY22 vs $0.47 FY21? -> Planned budget increase in FY22 (CR - No Actual Budget)

- **ROI – Actual SBU cost**
  - Number of SBUs delivered / Actual FY Budget (typically ~85% - 110M in FY21)
  - FY21 SBU == $0.47 vs. ROI == $0.42
- **Marginal SBU Cost**

  - SBU cost for last programs who added funding to HECC planed procurement FY21 (@75% utilization)
    - » $0.20 / SBU - 3 years
    - » $0.12 / SBU - 5 years
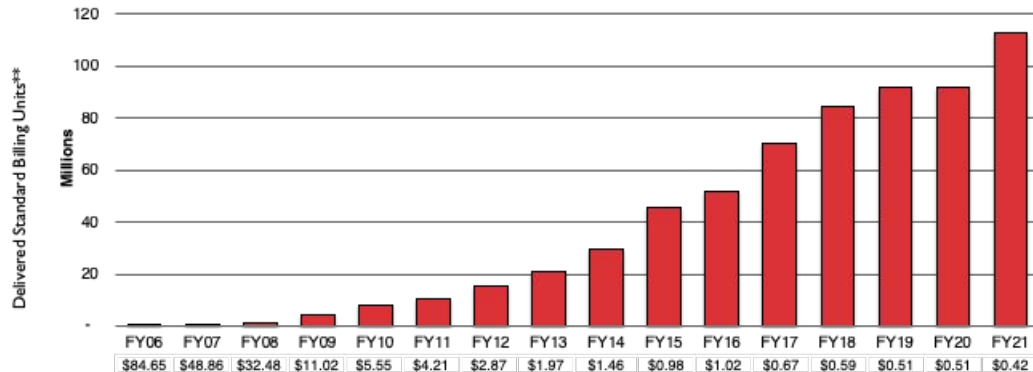    - » $0.09 / SBU - 7 years
- **Marginal ROI**
  - Likely - 10% better than marginal SBU cost
- **3-year payback time HECC vs cloud\***

  65 days on X nodes in cloud on-demand
  is equivalent to
  1095 days X nodes in HECC ( # of days in 3 years)
     (Substitute X for some number of nodes - e.g., 128)
  - HECC needs 6-12 months lead time
  - Based on recent CFD runs in cloud (AWS)
  - Comparable performance – no scaling advantage – 17x the cost vs. 3 year purchase at margin

Cost is predictable
Fixed Budget

Costs are all-inclusive: facilities, power, personnel, hardware, networking, maintenance, and storage

ROI – Actual SBU Cost



| | FY06 | FY07 | FY08 | FY09 | FY10 | FY11 | FY12 | FY13 | FY14 | FY15 | FY16 | FY17 | FY18 | FY19 | FY20 | FY21 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $84.65 | $48.86 | $32.48 | $11.02 | $5.55 | $4.21 | $2.87 | $1.97 | $1.46 | $0.98 | $1.02 | $0.67 | $0.59 | $0.51 | $0.51 | $0.42 |

# HECC Overview

NASA Ames' High-End Computing Capability (FY21 Budget: $47M)

- NASA's only cloud-scale private HPC cloud infrastructure
- Similar economics of scale to CSPs who run multiple hyperscale data center
- Robust power infrastructure from aero/windtunnel heritage (each >100MW)
- Current Pad/Power infrastructure can scale out to ~1 Billion SBUs (7.5x)
  - straight forward to double again
- > 620,000 CPU cores  > 614,000 GPU cores
- > 17,000 compute nodes - (delivered 110 million SBUs FY21)
- > 100 PB of on-line data storage
- > 350 PB of off-line tape data
- Supplementary analysis systems
- Scientific Consulting for Optimization and Help (significant code speed up)
- Data Analysis and Visualization – as a service



Traditional focus on modeling and simulation (Data Producer)
- Evolving Support for Hybrid Computing (On-prem/Cloud)
- Improvements around Latency (reservations and dedicated systems)

Growth in the size and nature of SMD data sets and SPD-41 motivating changes in HECC

- ECCO (Ocean Only Model Outputs) – 4 PB data set accessible through the NAS Data Portal
- GEOS Coupled – 3 PB of scratch space (https://gmao.gsfc.nasa.gov/GEOS_systems/)
- ECCO with GEOS5 – increased simulation output ~10 – 20 PB
  - Ocean Biology (e.g., predict whale migration/feeding patterns based on food sources estimated from global models)
- NASA's Earth Exchange (NEX) - 5.9 PB of on-line storage, used for data cache based on projects requirements for a given funding cycle.
- SBG is expected to collect 2.4 TB day$^{-1}$ and produce 40 TB day$^{-1}$ of data products